AD-785 031

# SITUATIONS, ACTIONS, AND CAUSAL LAWS

John McCarthy

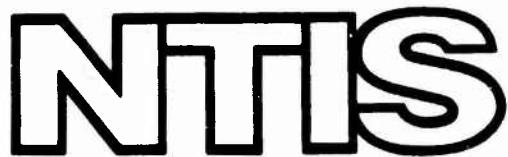Stanford University

Prepared for:

Advanced Research Projects Agency

3 July 1963

STANFORD ARTIFICIAL INTELLIGENCE PROJECT
MEMO NO.2

July 3, 1963

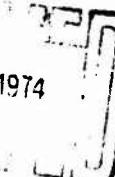## SITUATIONS, ACTIONS, AND CAUSAL LAWS

by John McCarthy

Abstract:  A formal theory is given concerning
situations, causality and the possibility
and effects of actions is given.  The
theory is intended to be used by the
Advice Taker, a computer program that is
to decide what to do by reasoning.  Some
simple examples are given of descriptions
of situations and deductions that certain
goals can be achieved.

## 1. INTRODUCTION

Although formalized theories have been devised to express the most important fields of mathematics and some progress has been made in formalizing certain empirical sciences, there is at present no formal theory in which one can express the kind of means-ends analysis used in ordinary life. The closest approach to such a theory of which I am aware is made by Freudenthal in Lincos [1].

Our approach to the artificial intelligence problem requires a formal theory. Namely, we believe that human intelligence depends essentially on the fact that we can represent in language facts about our situation, our goals, and the effects of the various actions we can perform. Moreover, we can draw conclusions from the facts to the effect that certain sequences of actions are likely to achieve our goals.

In Programs with Common Sense [2], I discussed the advantages of having a computer program, to be called the Advice Taker that would reason from collections of facts about its problem and derive statements about what it could do. The name Advice Taker came from the hope that its behavior could be improved by giving it advice in the form of new facts rather than by rewriting the program. The reader is referred to that paper for further information about the Advice Taker and to Minsky's paper Steps Towards Artificial Intelligence [3] for a general introduction to the subject.

The first requirement for the Advice Taker is a formal system in which facts about situations, goals and actions can be expressed and containing the general facts about means and ends as axioms. A start is made in this paper on providing a system meeting the following specifications

#1. General properties of causality and certain obvious but until now unformalized facts about the possibility and results of actions are given as axioms.

#2. It is a logical consequence of the facts of a situation and the general axioms that certain persons can achieve certain goals by taking certain actions.

#3. The formal descriptions of situations should correspond as closely as possible to what people may reasonably be presumed to know about them when deciding what to do.

## 2. SITUATIONS AND FLUENTS

One of the basic entities in our theory is the situation. Intuitively, a situation is the complete state of affairs at some instant of time. The laws of motion of a system determine from a situation all future situations. Thus a situation corresponds to the notion in physics of a point in phase space. In physics, laws are expressed in the form of differential equations which give the complete motion of the point in phase space.

Our system is not intended for the complete description of situations nor for the description of complete laws of motion. Instead, we deal with partial descriptions of situations and partial laws of motion. Moreover, the emphasis is on the simple qualitative laws of everyday life rather than on the quantitative laws of physics. As an example, take the fact that if it is raining and I go outside I will get wet.

Since a situation is a complete state of affairs we can never describe a situation completely, and therefore we provide no notation for doing so in our theory. Instead, we state facts about situations in the language of an extended predicate calculus. Examples of such facts are:

1.  raining (s)

    meaning that it is raining in situation s

2.  time (s) = 1963.7205

    giving the value of the time in situation s. It will usually prove convenient to regard the time as a function of the situation rather than vice versa. The reason for this is that the numerical value of the time is known and important only where the laws of physics are being used.

3.  at(I,home,s)  or  at(I,home)(s)

    meaning that I am at home in situation s. We prefer and will use the second of the given notations that isolates the situation variable since in most if not all cases we will be able to suppress it completely.

We shall not describe in this memorandum the logical system we intend to use. Basically, it is a predicate calculus, but we shall use the $\lambda$-notation and if necessary conditional expressions as in LISP or ALGOL. We shall extend the meaning of the Boolean operators to operate on predicates. Thus by

$$at(I,home) \wedge raining$$

we mean the same as

$$\lambda s.\, at(I,home)(s) \wedge raining(s)$$

A predicate or function whose argument is a situation will be called a _fluent_, the former being called a _propositional fluent_. Thus. raining, time, and at(I,home) are all fluents, the first and last being propositional fluents.

2

The term was used by Newton for a physical quantity that depends on time, and according to my limited understanding of what he meant, the present use of the term is justified.

In our formulas we will usually manage to use the fluents without explicitly writing variables representing situations. This corresponds to the use of random variables in probability theory without using variables representing points in the sample space even though random variables are supposed to be regarded as functions defined on a sample space.

In fact we shall go further and givenan interpretation of our theory as a sort of modal logic in which the fluents are not regarded ee functions et al!.

## 3. CAUSALITY

In order to express causal laws we introduce the second order predicate cause. The statement

$$\text{cause}(\pi)(s)$$

where $\pi$ is a propositional fluent is intended to mean that the situation s will lead in the future to a situation that satisfies the fluent $\pi$. Thus, cause($\pi$) is itself a propositional fluent. As an example of its use we write

$$\forall s.\forall p.\Big[\text{person}(p)\wedge\text{raining}\wedge\text{outside}(p)\supset\text{cause}(\text{wet}(p))\Big](s)$$

which asserts that a person who is outside when it is raining will get wet. We shall make the convention that if $\pi$ is a fluent then

$$\forall\pi$$

means the same as

$$\forall s.\,\pi(s).$$

With this convention we can write the previous statement as

$$\forall p. \text{person}(p)\wedge\text{raining}\wedge\text{outside}(p)\supset\text{cause}(\text{wet}(p)),$$

which suppresses explicit mention of situations. As a second example we give a special case of the law of falling bodies in the form:

$$\forall t.\forall b.\ \forall t'.\forall h\langle\text{real}(t)\wedge\text{real}(t')\wedge\text{real}(h)\wedge\text{body}(b)$$

$$\wedge\ \text{unsupported}(b)\wedge\Big[\text{height}(b)\ =\ h\Big]\wedge\Big[\tfrac{1}{2}gt^2<h\Big]\wedge$$

$$\Big[\text{time}\ =\ t'\Big]\supset\text{cause}(\text{height}(b)\ =\ h-\tfrac{1}{2}gt^2\wedge\text{time}=t'+t)$$

The concept of causality is intended to satisfy the two following general laws, which may be taken as axioms:-

C1.    $\forall.\ \text{cause}(\pi)\wedge\ [\forall.\pi\supset\rho]\supset\text{cause}(\rho)$

C2.    $\forall\ \text{cause}\ (\text{cause}(\pi))\supset\text{cause}(\pi)$

C3.    $\forall.\text{cause}(\pi_1)\vee\text{cause}(\pi_2)\supset\text{cause}(\pi_1\vee\pi_2)$

The fact that we can suppress explicit mention of situations has the following interesting consequence. Instead of regarding the $\pi$'s as predicates we may regard them as propositions and regard cause as a new modal operator. The operator $\forall$ seems then to be equivalent to the N (necessary) operator of ordinary modal logic.

Conversely, it would appear that modal logic of necessity might be regarded as a monadic predicate calculus where all quantifiers are over situations.

In the present case of causality, we seem to have our choice of how to proceed. Regarding the system as a modal logic seems to have the following two advantages.

4

1. If we use the predicate calculus interpretation we require second order predicate calculus in order to handle cause($\pi$)(s), while if we take the modal interpretation we can get by with first order predicate calculus.

2. We shall want decision procedures or at least proof procedures for as much of our system as possible. If we use the modal approach many problems will involve only substitution of constants for variables in universal statements and will therefore fall into a fairly readily decidable domain.

Another example of causality is given by a 2-bit binary counter that counts every second. In our formalism its behavior may be described by the statement:

$$\forall t\ \forall x_0 \forall x_1 /\ time = t \wedge bit00 = x_0 \wedge bit1 = x_1 \supset cause\ ($$

$$time = t+1 \wedge (bit\ 0 = x_0 \oplus 1) \wedge (bit\ 1 = x_1 \oplus (x_0 \wedge 1)))$$

In this example <u>time</u>, <u>bit00</u> and <u>bit11</u> are fluents while $t$, $x_0$ and $x_1$ are numerical variables. The distinction is made clearer if we use the more long-winded statement

$$\forall s \forall t \forall x_0 \forall x1.\ time(s) = t \wedge bit\ 0(s) = x_0 \wedge bit1(s) = x_1 \supset$$

$$cause(\lambda s'.\ time(s') = t+1 \wedge (bit0(s') = x_0 \oplus 1) \wedge (bit1(s') = x_1 \oplus (x_0 \wedge 1)))(s)$$

In this case however we can rewrite the statement in the form

$$\forall s.\ cause(\lambda s'. \left[ time(s') = time(s)+1 \right] \wedge \left[ bit0(s') = bit0(s) \oplus 1 \right] \wedge$$

$$\left[ bit1(s') = bit1(s) \oplus (bit\ 0(s) \wedge 1) \right])\ (s)$$

Thus we see that the suppression of explicit mention of the situations forced us to introduce the auxiliary quantities $t$, $x_0$ and $x_1$ which are required because we can no longer use functions of two different situations in the same formula. Nevertheless, the s-suppressed form may still be worthwhile since it admits the modal interpretation.

The time as a fluent satisfies certain axioms. The fact that there is only one situation corresponding to a given value of the time may be expressed by the axiom

T1. $\forall \pi \forall \rho \forall t.\ cause(time = t \wedge \pi) \wedge cause(time = t \wedge \rho) \supset$

cause $(time = t \wedge \pi \wedge \rho)$

Another axiom is

T2. $\forall t.\ real(t) \wedge t > time \supset cause(time = t)$

5

# 4. ACTIONS AND THE OPERATOR <u>can</u>

We shall regard the fact that a person performs a certain action in a situation as a propositional fluent. Thus

$$\text{moves (person, object, location) (s)}$$

is regarded as asserting that <u>person</u> <u>moves</u> <u>object</u> to <u>location</u> in the situation s. The effect of moving something is described by

$$\forall p \; \forall o \; \forall l. \; \text{moves } (p,o,l) \supset \text{cause } (at(o,l))$$

or in the long form

$$\forall s \; \forall p \; \forall o \; \forall l. \text{moves } (p,o,l)(s) \supset \text{cause } (\lambda s'.at \, (o,l)(s'))(s)$$

In order to discuss the ability of persons to achieve goals and to perform actions we introduce the operator <u>can</u>.

$$\text{can}(p,\pi) \; (s)$$

asserts that the person p can make the situation s satisfy. We see that can $(p,\pi)$ is a propositional fluent and that like <u>cause</u>, <u>can</u> may be regarded either as a second order predicate or a modal operator. Our most common use of <u>can</u> will be to assert that a person can perform a certain action. Thus we write

$$\text{can}(p, \text{moves } (p,o,l)) \; (s)$$

to assert that in situation s, the person p can move the object o to location l.

The operator <u>can</u> satisfies the axioms

K1. $\forall \pi \; \forall \rho \; \forall p. \; [\text{can } (p,\pi) \land (\pi \supset \rho) \supset \text{can}(p, \rho)$

K2. $\forall \pi \; \forall p_1 \; \forall p_2. \; [ \sim \text{can}(p_1, \pi) \land \text{can}(p_2, \sim \pi)]$

K3. $\forall p \; \forall \pi \forall \rho \; [\text{can}(p,\pi) \lor \text{can } (p, \rho) \supset \text{can}(p,\pi \lor \rho)]$

Using K1 and

$$\text{can}(p, \text{moves } (p, o, l))$$

and

$$\forall p \; \forall o \; \forall l. \text{ moves } (p, o, l) \supset \text{cause } (at(o,l))$$

we can deduce

$$\text{can } (p, \text{cause } (at \, (o,l)))$$

which shows that the operators <u>can</u> and <u>cause</u> often show up in the same formula.

The ability of people to perform joint actions can be expressed by formulas like

$$\text{can}(p_1, \text{can } (p_2, \text{marry } (p_1, p_2)))$$

which suggests the commutative axiom

K4. $\forall p_1 \; \forall p_2 \; \forall \pi. \; \text{can}(p_1, \text{can}(p_2, \pi)) \supset \text{can } (p_2, \text{can}(p_1, \pi))$

A kind of transitivity is expressed by the following:-

Theorem - From

1)   can(p, cause($\pi$))

and

2)   v. $\pi \supset$ can(p, cause($\rho$))

it follows that

3)   can(p , cause(can(p, cause($\rho$))))

Proof - Substitute can(p, cause($\rho$)) for $\rho$ in axiom C1 and substitute cause ($\pi$) for $\pi$ and cause(can(p, cause($\rho$))) for $\rho$ in axiom K1. The conclusion then follows by propositional calculus.

In order to discuss the achievement of goals requiring several consecutive actions we introduce canult(p,$\pi$) which is intended to mean that the person p can ultimately bring about a situation satisfying $\pi$. We connect it with can and cause by means of the axiom

KC1.    V. Vp V$\pi$. $\pi$V can(p,cause(canult(p,$\pi$)))$\supset$canult(p,$\pi$)

This axiom partially corresponds to the LISP-type recursive definition

canult(p,$\pi$)  =  $\pi$ V can(p, cause(canult(p,$\pi$)))

We also want the axiom

KC2.    VVpV$\pi$. cause(canult(p,$\pi$))$\supset$canult(p,$\pi$)

7

## 5. EXAMPLES

### 1. The Monkey can get the Bananas

The first example we shall consider is a situation in which a monkey is in a room where a bunch of bananas is hanging from the ceiling too high to reach. In the corner of the room is a box, and the solution to the monkey's problem is to move the box under the bananas and climb onto the box from which the bananas can be reached.

We want to describe the situation in such a way that it follows from our axioms and the description that the monkey can get the bananas. In this memorandum we shall not discuss the heuristic problem of how monkeys do or even might solve the problem. Specifically, we shall prove that

$$canult(monkey, has(monkey, bananas))$$

The situation is described in a very oversimplified way by the following seven statements:-

H1. $\forall u. \, place(u) \supset can(monkey, move(monkey, box, u))$

H2. $\forall u \, \forall v \, \forall p \, move \, (p,v,u) \supset cause(at(v,u))$

H3. $\forall \, can(monkey, climbs(monkey, box))$

H4. $\forall \, \forall u \forall v \forall p. \, at(v,u) \wedge climbs(p,v) \supset cause(at(v,u) \wedge on(p,v))$

H5. $\forall \, place(under(bananas))$

H6. $\forall \, at(box, under(bananas)) \wedge on(monkey, box) \supset can(monkey, reach(monkey, bananas))$

H7. $\forall \, \forall p \, \forall x. \, reach(p,x) \supset cause(has(p,x))$

The reasoning proceeds as follows: From 1 and 5 by substitution of under(bananas) for u and PC (propositional calculus) we get

1)     $can(monkey, move(box, under(bananas)))$

Using 1) and H2 and axiom C1, we get

2)     $can(monkey, cause(at(box, under(bananas))))$

Similarly H3 and H4 and C1 give

3)   $at(box) \, under(bananas)) \supset can(monkey, cause( \, at(box, under(bananas)) \wedge on(monkey, box)))$

Then H6 and H7 give

4)   $at(box, under(bananas)) \wedge on(monkey, box) \supset can(monkey, cause \, (has(monkey, bananas)))$

Now, Theorem 1 is used to combine 2), 3) and 4) to get

5)   $can(monkey, cause(can(monkey, cause(can(monkey, cause(has(monkey, bananas)))))))$

Using KC1, we reduce this to

$$canult(monkey, has(monkey, bananas)),$$

8

## 2. AN ENDGAME

A simple situation in a two person game can arise where player $p_1$ has two moves, but whichever he chooses player $p_2$ has a move that will beat him. This situation may be described as follows:-

1) $can(p_1, m_1) \wedge can(p_1, m_2) \wedge (m_1 \vee m_2)$

2) $\left[ m_1 \supset cause(\pi_1) \right] \wedge \left[ m_2 \supset cause(\pi_2) \right]$

3) $\forall . \pi_1 \vee \pi_2 \supset \left[ can(p_2, n_1) \wedge can(p_2, n_2) \wedge (n_1 \vee n_2) \right]$

4) $\forall . (\pi_1 \wedge n_1) \vee (\pi_2 \wedge n_2) \supset cause(win(p_2))$

We would like to be able to draw the conclusion

3) $canult(p_2, win(p_2))$

We proceed as follows: From 1) and 2) we get

4) $cause(\pi_1) \vee cause(\pi_2)$

and we use axiom C3 to get

5) $cause(\pi_1 \vee \pi_2)$

Next we weaken 3) to get

6) $\forall . \pi_1 \supset can(p_2, n_1)$  and

7) $\forall . \pi_2 \supset can(p_2, n_2)$

and then we use K1 to get

8) $\forall . \pi_1 \supset can(p_2, \pi_1 \wedge n_1)$  and

9) $\forall . \pi_2 \supset can(p_2, \pi_2 \wedge n_2)$

The propositional calculus gives

10) $\forall . \pi_1 \vee \pi_2 \supset can(p_2, \pi_1 \wedge n_1)$   $can(p_2, \pi_2 \wedge n_2)$

and using K3 we get

11) $\forall . \pi_1 \vee \pi_2 \supset can(p_2, (\pi_1 \wedge n_1) \vee (\pi_2 \wedge n_2))$

which together with 4) and K1 gives

12) $\forall . \pi_1 \vee \pi_2 \supset can(p_2, cause(win(p_2)))$

which together with 5) and C1 gives

13) $cause(can(p_2, cause(win(p_2))))$

Using the axioms for _canult_ we now get

14) $canult(p_2, win(p_2))$.

9

## 6. NOTE

After finishing the bulk of this memorandum I came across <u>The Syntax of Time Distinctions</u>  4  by A.N.Prior.  Prior defines model operators P and F where

P($\mathcal{T}$)  means  'it has been the case that $\mathcal{T}$' end,

F($\mathcal{T}$)  means  'it will be the case that $\mathcal{T}$'

He subjects these operators to a number of axioms and rules of inference in close analogy to the well-known [5] modal logic of possibility.  He also interprets this logic in a restricted predicate calculus where the variables range over times. He then extends his logic to include a somewhat undetermined future and claims (unconvincingly) that this logic cannot be interpreted in predicate calculus.

I have not yet made a detailed comparison of our logic with Prior's, but here are some tentative conclusions.

1.  The causality logic should be extended to allow inference about the past.

2.  Causality logic should be extended to allow inference that certain propositional fluents will always hold.

3.  cause($\mathcal{T}$) satisfies the axioms for his F($\mathcal{T}$) which means that his futurity theory possesses, from his point of view, non-standard models.  Namely, a collection of functions $p_1(t)$,$p_2(t)$ may satisfy his futurity axioms and assign truth to $p(1) \wedge \sim(Fp)(o)$.  In our system this is okay because something can happen without being caused to happen.

4.  If we combine his past and futurity axioms, our system will no longer fit his axioms and

PF1.     $p \supset \sim F(\sim P(p))$

PF2.     $p \supset \sim P(\sim F(p))$

since we do not wish to say that whatever is, was always inevitable.

# REFERENCES

1. Freudenthal, H.A., "LINCOS, Design of a Language for Cosmic Intercourse Part I," Amsterdam (1960)

2. McCarthy, J. "Programs with Common Sense", Proceedings Symposium on Mechanization of Thought Processes, Her Majesty's Stationery Office, London, England; 1959

3. Minsky, M.L., "Steps Toward Artificial Intelligence", Proceedings of the IRE, Vol. 29, No. 1, January 1961

4. Prior, A.N "The Syntax of Time Distinctions" Franciscan Studies (1958)

5. von Wright, G.H. "An Essay in Modal Logic", Amsterdam (1951)